

**IN THE UNITED STATES PATENT AND TRADEMARK OFFICE**

Applicant: Tadashi EMORI, et al.  
Title: SPEAKER'S VOICE  
RECOGNITION SYSTEM,  
METHOD AND RECORDING  
MEDIUM  
Appl. No.: Unassigned  
Filing Date: October 25, 2000  
Examiner: Unassigned  
Art Unit: Unassigned



**CLAIM FOR CONVENTION PRIORITY**

Assistant Commissioner for Patents  
Washington, D.C. 20231

Sir:

The benefit of the filing date of the following prior foreign application filed in the following foreign country is hereby requested, and the right of priority provided in 35 U.S.C. § 119 is hereby claimed.

In support of this claim, filed herewith is a certified copy of said original foreign application:

Japanese Patent Application  
No. 11-304685 filed 26 OCTOBER 1999.

Respectfully submitted,

Date: October 25, 2000

FOLEY & LARDNER  
Washington Harbour  
3000 K Street, N.W., Suite 500  
Washington, D.C. 20007-5109  
Telephone: (202) 672-5407  
Facsimile: (202) 672-5399

By Phillips J. Artisola Reg. No. 38,819  
for / David A. Blumenthal  
Attorney for Applicant  
Registration No. 26,257

日本国特許庁

PATENT OFFICE  
JAPANESE GOVERNMENT

03  
JCS20 U.S. PRO  
09/695067  
10/25/00

別紙添付の書類に記載されている事項は下記の出願書類に記載されて  
いる事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed  
with this Office.

出願年月日  
Date of Application:

1999年10月26日

出願番号  
Application Number:

平成11年特許願第304685号

願人  
Applicant(s):

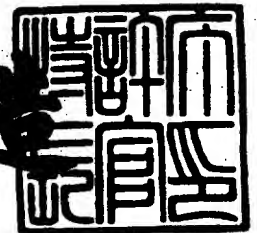
日本電気株式会社

CERTIFIED COPY OF  
PRIORITY DOCUMENT

2000年 7月14日

特許庁長官  
Commissioner,  
Patent Office

及川耕造



出証番号 出証特2000-3054464

【書類名】 特許願

【整理番号】 33509635

【提出日】 平成11年10月26日

【あて先】 特許庁長官 殿

【国際特許分類】 G10L 3/00  
G10L 3/02

【発明者】

    【住所又は居所】 東京都港区芝五丁目 7 番 1 号 日本電気株式会社内

    【氏名】 江森 正

【発明者】

    【住所又は居所】 東京都港区芝五丁目 7 番 1 号 日本電気株式会社内

    【氏名】 篠田 浩一

【特許出願人】

    【識別番号】 000004237

    【氏名又は名称】 日本電気株式会社

【代理人】

    【識別番号】 100080816

    【弁理士】

    【氏名又は名称】 加藤 朝道

    【電話番号】 045-476-1131

【手数料の表示】

    【予納台帳番号】 030362

    【納付金額】 21,000円

【提出物件の目録】

    【物件名】 明細書 1

    【物件名】 図面 1

    【物件名】 要約書 1

    【包括委任状番号】 9304371

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 音声認識装置及び方法ならびに記録媒体

【特許請求の範囲】

【請求項 1】

音声信号のスペクトルを周波数軸上で伸縮を行うスペクトル変換部が、  
入力音声信号をケプストラムを含む入力パターンに変換する分析部と、  
標準パターンを記憶する標準パターン記憶部と、  
前記入力パターンと前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定部と、  
前記伸縮パラメータを用いて前記入力パターンを変換する変換部と、  
を備えたことを特徴とする音声認識装置。

【請求項 2】

入力音声信号をケプストラムを含む入力パターンに変換する分析部と、  
標準パターンを記憶する標準パターン記憶部と、  
前記入力パターンと前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定部と、  
前記伸縮パラメータを用いて前記入力パターンを変換する変換部と、  
前記変換部から出力された伸縮後の前記入力パターンと前記標準パターンとの距離を算出し最も距離の小さい前記標準パターンを認識結果として出力する整合部と、  
を備えたことを特徴とする音声認識装置。

【請求項 3】

前記変換部が、スペクトルの伸縮の形態を定めるワーピング関数によるスペクトルの周波数軸上の伸縮をケプストラム空間で行うことで、スペクトルの伸縮を行う、ことを特徴とする請求項 1 又は 2 記載の音声認識装置。

【請求項 4】

前記伸縮推定部が、スペクトルの伸縮の形態を定めるワーピング関数によって定められたスペクトルの周波数軸上の伸縮を、ケプストラム空間における HMM（隠れマルコフモデル）の最尤推定から導出される推定を用いて伸縮パラメータ

の推定を行う、ことを特徴とする請求項 1 乃至 3 のいずれか一に記載の音声認識装置。

【請求項 5】

学習用音声記憶しておく学習音声記憶部と、  
前記学習音声記憶部から学習用音声信号を入力し前記学習用音声信号をケプストラムを含む入力パターンに変換する分析部と、  
標準パターンを記憶する標準パターン記憶部と、  
前記入力パターンと前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定部と、  
前記伸縮パラメータを用いて入力パターンを変換する変換部と、  
前記学習用音声に対し、前記変換部が出力した伸縮後の入力パターンと前記標準パターンとを用いて、前記標準パターン記憶部の標準パターンを更新する標準パターン推定部と、  
前記変換部から出力される伸縮後の入力パターンと前記標準パターンとを用いて、距離を算出し、距離の変化量を監視する尤度判定部と、  
を備えたことを特徴とする標準パターン学習装置。

【請求項 6】

前記変換部が、スペクトルの伸縮の形態を定めるワーピング関数によるスペクトルの周波数軸上の伸縮をケプストラム空間で行うことでスペクトルの伸縮を行う、ことを特徴とする請求項 5 記載の標準パターン学習装置。

【請求項 7】

前記伸縮推定部が、スペクトルの伸縮の形態を定めるワーピング関数によって定められたスペクトルの周波数軸上の伸縮をケプストラム空間での HMM の最尤推定から導出された推定方法を用いて伸縮パラメータの推定を行うことを特徴とする請求項 5 又は 6 記載の標準パターン学習装置。

【請求項 8】

入力音声信号をケプストラムを含む入力パターンに変換する分析部と、  
標準パターンを記憶する標準パターン記憶部と、  
前記入力パターンと前記標準パターンとを用いて周波数軸方向の伸縮パラメー

タを出力する伸縮推定部と、

前記伸縮パラメータを用いて入力パターンを変換する変換部と、

前記変換部から出力された伸縮後の前記入力パターンの時系列を逆変換して時間領域波形を出力する逆変換部と、

を備えたことを特徴とする声質変換装置。

【請求項 9】

音声信号のスペクトルを周波数軸上で伸縮を行うスペクトル変換部を構成するコンピュータにおいて、

(a) 入力音声信号をケプストラムを含む入力パターンに変換する分析処理と

(b) 前記入力パターンと、標準パターンを記憶する標準パターン記憶部に記憶されている前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定処理と、

(c) 前記伸縮パラメータを用いて前記入力パターンを変換する変換処理と、  
の前記 (a) 乃至 (c) の処理を前記コンピュータで実行させるためのプログラムを記録する記録媒体。

【請求項 10】

音声信号のスペクトルを周波数軸上で伸縮を行い音声認識を行う装置を構成するコンピュータにおいて、

(a) 入力音声信号をケプストラムを含む入力パターンに変換する分析処理と

(b) 前記入力パターンと、標準パターンを記憶する標準パターン記憶部に記憶されている前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定処理と、

(c) 前記伸縮パラメータを用いて前記入力パターンを変換する変換処理と、

(d) 前記出力された伸縮後の前記入力パターンと前記標準パターンとの距離を算出し最も距離の小さい前記標準パターンを認識結果として出力する整合処理と、

の前記 (a) 乃至 (d) の処理を前記コンピュータで実行させるためのプログ

ラムを記録した記録媒体。

【請求項 11】

請求項 10 記載の記録媒体において、前記変換処理が、スペクトルの伸縮の形態を定めるワーピング関数によるスペクトルの周波数軸上の伸縮をケプストラム空間で行うことで、スペクトルの伸縮を行う処理を前記コンピュータで実行させるためのプログラムを記録した記録媒体。

【請求項 12】

請求項 10 記載の記録媒体において、前記伸縮推定処理が、スペクトルの伸縮の形態を定めるワーピング関数によって定められたスペクトルの周波数軸上の伸縮を、ケプストラム空間における HMM（隠れマルコフモデル）の最尤推定から導出される推定を用いて伸縮パラメータの推定を行う、処理を前記コンピュータで実行させるためのプログラムを記録した記録媒体。

【請求項 13】

学習用音声から標準パターンを学習する装置を構成するコンピュータにおいて

（a）学習用音声を記憶しておく学習音声記憶部に記憶されている学習用音声を入力しケプストラムを含む入力パターンに変換する分析処理と、

（b）前記入力パターンと、標準パターンを記憶する標準パターン記憶部に記憶されている前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定処理と、

（c）前記伸縮パラメータを用いて前記入力パターンを変換する変換処理と、

（d）学習用音声に対し、変換処理が出力した伸縮後の入力パターンと、標準パターンを利用して、前記標準パターンを更新する標準パターン推定処理と、

（e）前記変換処理で変換された伸縮後の入力パターンと前記標準パターンを利用して、距離を算出し、距離の変化量を監視する尤度判定処理と、

の前記（a）乃至（e）の処理を前記コンピュータで実行させるためのプログラムを記録する記録媒体。

【請求項 14】

請求項 1 3 記載の記録媒体において、前記変換処理が、スペクトルの伸縮の形態を定めるワーピング関数によるスペクトルの周波数軸上の伸縮をケプストラム空間で行うことで、スペクトルの伸縮を行う処理を前記コンピュータで実行させるためのプログラムを記録した記録媒体。

【請求項 1 5】

請求項 1 3 記載の記録媒体において、前記伸縮推定処理が、スペクトルの伸縮の形態を定めるワーピング関数によって定められたスペクトルの周波数軸上の伸縮を、ケプストラム空間における HMM（隠れマルコフモデル）の最尤推定から導出される推定を用いて伸縮パラメータの推定を行う、処理を前記コンピュータで実行させるためのプログラムを記録した記録媒体。

【請求項 1 6】

音声信号のスペクトルを周波数軸上で伸縮を行う装置を構成するコンピュータにおいて、

- (a) 入力音声信号をケプストラムを含む入力パターンに変換する分析処理と、
- (b) 前記入力パターンと、標準パターンを記憶する標準パターン記憶部に記憶されている前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定処理と、
- (c) 前記伸縮パラメータを用いて前記入力パターンを変換する変換処理と、
- (d) 前記出力された伸縮後の前記入力パターンの時系列を逆変換して時間領域波形を出力する逆変換処理と、

の前記 (a) 乃至 (d) の処理を前記コンピュータで実行させるためのプログラムを記録する記録媒体。

【請求項 1 7】

音声信号のスペクトルを周波数軸上で伸縮を行うスペクトル変換を行う音声認識方法において、

入力音声信号をケプストラムを含む入力パターンに変換する第 1 のステップと

前記入力パターンと、標準パターンを記憶する標準パターン記憶部に記憶され



ている前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する第2のステップと、

前記伸縮パラメータを用いて前記入力パターンを変換する第3のステップと、  
を含むことを特徴とするスペクトル変換方法。

【請求項18】

入力音声信号をケプストラムを含む入力パターンに変換する第1のステップと

、  
前記入力パターンと、標準パターンを記憶する標準パターン記憶部に記憶されている前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する第2のステップと、

前記伸縮パラメータを用いて前記入力パターンを変換する第3のステップと、  
前記出力された伸縮後の前記入力パターンと前記標準パターンとの距離を算出し最も距離の小さい前記標準パターンを認識結果として出力する第4のステップと、

を含むことを特徴とする音声認識方法。

【請求項19】

前記第3のステップにおいて、スペクトルの伸縮の形態を定めるワーピング関数によるスペクトルの周波数軸上の伸縮をケプストラム空間で行うことで、スペクトルの伸縮を行う、ことを特徴とする請求項17又は18記載の音声認識方法。

【請求項20】

前記第2のステップにおいて、スペクトルの伸縮の形態を定めるワーピング関数によって定められたスペクトルの周波数軸上の伸縮を、ケプストラム空間におけるHMM（隠れマルコフモデル）の最尤推定から導出される推定を用いて伸縮パラメータの推定を行う、ことを特徴とする請求項17乃至19のいずれかに記載の音声認識方法。

【請求項21】

学習用音声を記憶しておく学習音声記憶部に記憶されている学習用音声を入力しケプストラムを含む入力パターンに変換する第1のステップと、

前記入力パターンと、標準パターンを記憶する標準パターン記憶部に記憶されている前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する第 2 のステップと、

前記伸縮パラメータを用いて前記入力パターンを変換する第 3 のステップと、

前記学習用音声に対し、変換処理が出力した伸縮後の入力パターンと、標準パターンを利用して、前記標準パターンに更新する第 4 のステップと、

伸縮後の入力パターンと標準パターンを利用して、距離を算出し、距離の変化量を監視し尤度判定処理を行う第 5 のステップと、

を含むことを特徴とする標準パターン学習方法。

#### 【請求項 2 2】

前記第 3 のステップにおいて、スペクトルの伸縮の形態を定めるワーピング関数によるスペクトルの周波数軸上の伸縮をケプストラム空間で行う、ことを特徴とする請求項 2 1 記載の標準パターン学習方法。

#### 【請求項 2 3】

前記第 2 のステップにおいて、スペクトルの伸縮の形態を定めるワーピング関数によって定められたスペクトルの周波数軸上の伸縮を、ケプストラム空間における HMM（隠れマルコフモデル）の最尤推定から導出される推定を用いて伸縮パラメータの推定を行う、ことを特徴とする請求項 2 1 記載の標準パターン学習方法。

#### 【請求項 2 4】

音声信号のスペクトルを周波数軸上で伸縮を行うスペクトル変換を行う音声認識方法において、

ワーピング関数で定義された、入力音声信号のスペクトルの伸縮を、ケプストラム上で行い、スペクトルの周波数軸の伸縮の度合いを、前記ワーピング関数に含まれる伸縮パラメータで決定し、伸縮パラメータの値を、話者毎に、最適な値を決定する、ことを特徴とする音声認識方法。

#### 【発明の詳細な説明】

#### 【0 0 0 1】

#### 【発明の属する技術分野】

本発明は、不特定話者の音声認識装置、音声認識方法と音響モデルの学習方法および音声認識用プログラムを記録した記録媒体に関し、特に、周波数軸上の話者性を正規化できる音声認識装置および正規化学習装置、音声認識方法および正規化学習方法、音声認識用プログラムおよび正規化学習プログラムを記録した記録媒体に関する。

#### 【0002】

##### 【従来の技術】

従来の音声認識装置のスペクトル変換部として、例えば特開平6-214596号公報（「文献1」という）や、1997年にアイ・イー・イー・イー（IEEE）から刊行された「インターナショナル・カンファレンス・オン・アコースティックス・スピーチ・アンド・シグナル・プロシーディング」の論文集の第1039頁から第1042頁に掲載された「スピーカー・ノーマリゼーション・ベースド・オン・フリークエンシー・ワーピング」と題するザン達による論文（Puming Zhan and Martin Westphal "Speaker Normalization Based On Frequency Warping", ICASSP, 1039-1042, 1997）（「文献2」という）等の記載が参照される。

#### 【0003】

例えば、上記文献1には、予め定められた複数の異なる周波数特性補正係数に基づいて、入力された音声信号の周波数特性を補正する周波数補正手段と、予め定められた複数の周波数軸変換計数に基づいて、入力された音声信号の周波数軸を変換する周波数軸変換手段と、入力された音声信号の特徴量を入力音声特徴量として抽出する特徴量抽出手段と、標準音声特徴量を保持している標準音声記憶手段と、周波数特性補正手段、周波数軸変換手段、特徴量抽出手段により処理され得られた入力音声特徴量と標準音声記憶手段に保持されている標準音声特徴量とを照合する照合手段と、を備え、話者適応フェーズと音声認識フェーズの機能を有する音声認識装置であって、話者適応フェーズにおいて、既知なる発生内容の未知なる話者の入力音声信号に対して、複数の異なる周波数特性補正係数、複数の異なる周波数軸変換係数の各々の係数ごとに、周波数特性補正手段、周波数軸変換手段、特徴量抽出手段の処理を行わせて、各々の係数毎の入力音声特徴量

と既知なる発生内容と同一内容の標準音声特徴量と照合し、各々の係数のうちから、最小距離を与える 1 つの周波数特性補正係数と 1 つの周波数軸変換係数を選択し、音声認識フェーズでは、選択された周波数特性補正係数と周波数軸変換係数を用いて、入力音声特徴量を求め、これを標準音声特徴量と照合し音声認識処理を行う構成が提案されている。

## 【 0 0 0 4 】

これら従来の音声認識装置において、スペクトル変換部は、話者の性別や年齢や体格などの個人差に対して、音声信号のスペクトルを周波数軸上で伸縮を行うことにより、認識性能の改善を行う。

## 【 0 0 0 5 】

そしてスペクトル変換部では、周波数軸上で伸縮を行うために、適当なパラメータでその伸縮の概形を変えることのできる関数を定義し、音声信号のスペクトルを、ワーピング関数を用いて周波数軸上の伸縮を行う。

## 【 0 0 0 6 】

周波数軸上での伸縮を行うための関数を「ワーピング関数」という。また、ワーピング関数の概形を定義するパラメータを「伸縮パラメータ」という。

## 【 0 0 0 7 】

従来、ワーピング関数の伸縮パラメータ（「ワーピングパラメータ」と略記する）として、複数個のワーピングパラメータの値を用意し、1 つ 1 つの値を用いて音声信号のスペクトルを周波数軸上で伸縮を行い、伸縮されたスペクトルを用いて入力パターンを計算し、入力パターンと標準パターンを用いて距離を求め、該距離が最小になる値を、認識時のワーピングパラメータの値としている。

## 【 0 0 0 8 】

従来の音声認識装置のスペクトル変換部について、図面を参照して以下に説明する。図 9 は、従来の音声認識装置のスペクトル変換部の構成の一例を説明するための図である。図 9 を参照すると、この従来のスペクトル変換部は、FFT（Fast Fourier Transform；高速フーリエ変換）部 3 0 1 と、伸縮パラメータ記憶部 3 0 2 と、周波数変換部 3 0 3 と、入力パターン計算部 3 0 4 と、整合部 3 0 6 と、標準パターン 3 0 5 と、伸縮パラメータ選択部 3 0 7 と、を備えている。

【0 0 0 9】

F F T 部 3 0 1 は、入力された音声信号を一定の時間毎に切り出してフーリエ変換を行い、周波数スペクトルを求める。

【0 0 1 0】

伸縮パラメータ記憶部 3 0 2 は、周波数の伸縮を決定する伸縮パラメータの値を複数記憶している。

【0 0 1 1】

周波数変換部 3 0 3 は、伸縮パラメータによりその概形が決定されたワーピング関数を用い、F F T 部 3 0 1 から出力されるスペクトルに対して、周波数の伸縮を行い、周波数の伸縮が行われたスペクトルを伸縮スペクトルとして出力する。

【0 0 1 2】

入力パターン計算部 3 0 4 は、周波数変換部 3 0 3 から出力された伸縮スペクトルを用い、入力パターンを計算し、出力する。

【0 0 1 3】

入力パターンは、ケプストラムなど音響的特徴を表すパラメータの時系列などを表す。

【0 0 1 4】

標準パターンは、入力パターンなどを数多く用い、同じクラスに属する音素単位の入力パターンを、ある種の平均化の手法を用い作成されるものである。標準パターンの作成については、1 9 9 5 年に N T T アドバンステクノロジー株式会社から刊行された「古井監訳、音声認識の基礎」（上）の第 6 3 頁（「文献 3」という）の記載が参照される。

【0 0 1 5】

標準パターンの種類には、認識アルゴリズムにより、例えば D P (Dynamic Programming; ダイナミックプログラミング) マッチングの場合、音素の時系列順に入力パターン並べられた時系列標準パターン、HMM (Hidden Markov Model; 隠れマルコフモデル) の場合、状態の系列やその接合情報となる。

【0 0 1 6】

整合部306は、FFT部301に入力された発声の内容に合わせた標準パターン305と入力パターンを用い、距離を求める。算出される距離は、標準パターンがHMMの場合、尤度に相当し、DPマッチングの場合、最適経路の距離に相当する。

【0017】

伸縮パラメータ選択部307は、整合部306で求められた整合性から最も整合性のとれている伸縮パラメータを選択する。

【0018】

図10は、従来のスペクトル整合部の処理を説明するための流れ図である。図0及び図10を参照して、従来のスペクトル変換部の動作について説明する。

【0019】

FFT部301は、音声信号をFFT演算することにより、スペクトルを求め出力する（図10のステップD101）。

【0020】

周波数変換部303は、入力された伸縮パラメータ（D106）を用いてスペクトルの周波数軸の伸縮を行う（ステップD102）。

【0021】

入力パターン計算部304は、周波数軸の伸縮されたスペクトルを用い、入力パターンを計算する（ステップD103）。

【0022】

整合部305は、標準パターン（D107）と入力パターンの距離を求める（D104）。

【0023】

ステップD101からD104までの一連の処理を、伸縮パラメータ記憶部302に記憶されている全ての伸縮パラメータの値について行う（ステップD105）。

【0024】

すなわち、伸縮パラメータ記憶部302に、伸縮パラメータの値として10個の値が記憶されている場合、ステップD101からD104の処理を10回繰り返す。

返し、10通りの距離を求める。

【0025】

伸縮パラメータ選択部307は、全ての伸縮パラメータに対応する距離を比較し、最も距離の小さい場合の、伸縮パラメータを選択する（ステップD108）

【0026】

【発明が解決しようとする課題】

しかしながら、上記した従来のスペクトル変換部は、下記記載の問題点を有している。

【0027】

第1の問題点は、伸縮パラメータの値を決定する際に、演算量が大きくなる、ということである。

【0028】

その理由は、従来のスペクトル変換部においては、伸縮パラメータの値を複数個用意しておき、その数の分だけ、

- ・FFT、
- ・スペクトルの周波数の伸縮、
- ・入力パターンの計算と、
- ・距離を計算する処理

などを繰り返し行う必要があるためである。

【0029】

第2の問題点は、周波数の伸縮による認識装置などへの効果が、十分得られない可能性が有る、ということである。

【0030】

その理由は、伸縮パラメータの値は、予め定められた値しか持たないので、未知の話者に対して、それらの値の中に最適なものがあるとは限らないためである

【0031】

したがって、本発明は、上記問題点に鑑みてなされたものであって、その主た

る目的は、少ない演算量で話者毎に最適な伸縮パラメータの値を計算可能とし、性能を向上する音声認識装置及び方法並びに記録媒体を提供することにある。これ以外の本発明の目的、特徴等は以下の説明で直ちに明らかとされるであろう。

#### 【0032】

##### 【課題を解決するための手段】

前記目的を達成する本発明は、入力音声信号をケプストラムを含む入力パターンに変換する分析部と、標準パターンを記憶する標準パターン記憶部と、前記入力パターンと前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定部と、前記伸縮パラメータを用いて前記入力パターンを変換する変換部と、前記変換部が出力した伸縮後の前記入力パターンと前記標準パターンとの距離を算出し、最も距離の小さい前記標準パターンを認識結果として出力する整合部と、を備えている。

#### 【0033】

本発明は、学習用音声記憶しておく学習音声記憶部と、前記学習音声記憶部から学習用音声信号を入力し前記学習用音声信号をケプストラムを含む入力パターンに変換する分析部と、標準パターンを記憶する標準パターン記憶部と、前記入力パターンと前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定部と、前記伸縮パラメータを用いて入力パターンを変換する変換部と、前記標準パターンを記憶する標準パターン記憶部と、学習用音声に対し、変換部が出力した伸縮後の入力パターンと標準パターンを利用して、標準パターンに更新する標準パターン推定部と、伸縮後の入力パターンと標準パターンを利用して、距離を算出し、距離の変化量を監視する尤度判定部と、を備える。

#### 【0034】

本発明は、入力音声信号をケプストラムを含む入力パターンに変換する第1の工程と、前記入力パターンと、標準パターンを記憶する標準パターン記憶部に記憶されている前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する第2の工程と、前記伸縮パラメータを用いて前記入力パターンを変換する第3の工程と、前記出力された伸縮後の前記入力パターンと前記標準パターンとの距離を算出し最も距離の小さい前記標準パターンを認識結果として出力する第4



の工程と、を含む。

【 0 0 3 5 】

【発明の実施の形態】

本発明の実施の形態について説明する。本発明は、入力音声信号をケプストラムを含む入力パターンに変換する分析部（１）と、標準パターンを記憶する標準パターン記憶部（４）と、入力パターンと標準パターンを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定部（３）と、前記伸縮パラメータを用いて入力パターンを変換する変換部（２）と、を備える。

【 0 0 3 6 】

また本発明は、変換部（２）によって変換された前記入力パターンと標準パターンとの距離を計算し、最も距離の小さい標準パターンを認識結果として出力する整合部（認識部 1 0 1）を備える。

【 0 0 3 7 】

伸縮推定部（３）は、入力パターンに含まれるケプストラムを用いて、伸縮パラメータを推定する。このため、本発明によれば、伸縮パラメータを決定する際に、前もって様々な値を記憶保持しておく必要も無く、また様々な値について距離計算を行う必要もない。

【 0 0 3 8 】

また本発明の装置は、学習用音声記憶部（２０１）と、学習音声記憶部（２０１）から学習用音声を入力しケプストラムを含む入力パターンに変換する分析部（１）と、標準パターンを記憶する標準パターン記憶部（４）と、前記入力パターンと前記標準パターンを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定部（３）と、前記伸縮パラメータを用いて入力パターンを変換する変換部（２）と、前記標準パターンを記憶する標準パターン記憶部と、学習用音声に対し、変換部が出力した伸縮後の入力パターンと標準パターンを利用して、標準パターンに更新する標準パターン推定部（２０２）と、伸縮後の入力パターンと標準パターンを利用して、距離を算出し、距離の変化量を監視する尤度判定部（２０３）と、を備える。

【 0 0 3 9 】

## 【実施例】

次に、本発明の実施例について図面を参照して詳細に説明する。図 1 は、本発明の第 1 の実施例の音声認識装置のスペクトル変換部の構成を示す図である。図 1 を参照すると、本発明の第 1 の実施例をなす音声認識装置のスペクトル変換部は、分析部 1 と、変換部 2 と、伸縮推定部 3 と、標準パターン記憶部 4 とを備えて構成されている。

## 【0040】

分析部 1 は、一定の時間毎に音声信号を切り出し、FFT（高速フーリエ変換）や LPC（linear predictive coding; 線形予測符号）分析などを用い、スペクトル成分を求め、人間の聴覚を考慮したメルスケールに変換し、そのメルスペクトルの成分の包絡成分を抽出するためのメルケプストラムを求め、そのメルケプストラムや、その変化量、変化量の変化量などを、入力パターンとして出力する。

## 【0041】

変換部 2 は、周波数の伸縮を、入力パターン中のメルケプストラムを変換することにより行う。ここで、変換部 2 で行われる変換の一例について詳細に説明する。

## 【0042】

次式（1）で示される 1 次の全域通過型のフィルタによる周波数の変換は、例えば、1972 年にアイ・イー・イー・イー（IEEE）から発行された「プロシーディングス・オブ・アイ・イー・イー・イー」のボリューム 60 の第 6 号の第 681 頁から 619 頁に掲載された「ディスクリート・リプレゼンテーション・オブ・シグナルス」と題するオッペンハイムによる論文（Oppenheim : "Discrete Representation of Signals," Proc. IEEE, 60, 681-691, June 1972）（「文献 4」という）によると、ケプストラム（記号  $c$ 、添え数字はケプストラムの次元数とする）を用いた次式（2）で示される再帰的表現で表すことができる。

【 0 0 4 3 】

$$\hat{z}^{-1} = \frac{z^{-1} - \alpha}{1 - \alpha z^{-1}} \quad \cdots(1)$$

$$\hat{c}_m^{(i)} = \begin{cases} c_{-i} + \alpha \hat{c}_0^{(i-1)}, & m = 0 \\ (1 - \alpha^2) \hat{c}_{m-1}^{(i-1)} + \alpha \hat{c}_1^{(i-1)}, & m = 1 \\ \hat{c}_{m-1}^{(i-1)} + \alpha (\hat{c}_m^{(i-1)} - \hat{c}_{m-1}^{(i)}), & m \geq 2 \end{cases}$$

$$i = -\infty, \dots, -1, 0. \quad \cdots(2)$$

【 0 0 4 4 】

式 (2) よるケプストラム空間での変換が、上式 (1) によるスペクトルの周波数の変換と等価になる。

【 0 0 4 5 】

そのため、変換部 1 0 2 は、周波数の伸縮のために直接スペクトルを用いるのではなく、式 (1) をワーピング関数とし、上式 (1) 中の  $\alpha$  を伸縮パラメータとし、式 (1) から導かれる式 (2) による変換を、入力パターンに対して施すことにより、スペクトルの周波数の伸縮を行う。

【 0 0 4 6 】

変換された入力パターンは、変換入力パターンとして出力される。

【 0 0 4 7 】

標準パターン記憶部 4 には、標準パターンが記憶される。標準パターンは、単語や音素などを単位とした音韻情報として、隠れマルコフモデル (「HMM」という) や、音素の時系列などの時系列標準パターンなどでも代用できる。本実施例では、標準パターンは、隠れマルコフモデルとする。HMMを構成する情報は、連続ガウス分布における平均ベクトルや分散、さらに状態と状態の遷移確率等のなどがある。

【 0 0 4 8 】

伸縮推定部 (「伸縮パラメータ推定部」ともいう) 3 は、分析部 1 に入力された音声信号に対応する HMM を用い、入力パターンとのアライメントを求める。ここで、アライメントとは、HMM の各時刻と各状態における事後確率のことである。

【0049】

アライメントの求め方としては、例えば、1995年にNTTアドバンステクノロジー株式会社から出版された「古井監訳、音声認識の基礎（下）」の第102頁乃至第185頁（「文献5」という）に記載されている、ビタビアルゴリズムやフォワードバックワードアルゴリズム等の公知の手法を用いて求めることができる。

【0050】

求められたアライメントとHMMと入力パターンを用い、伸縮パラメータを計算する。伸縮パラメータの計算は、式（4）によって求められる。

【0051】

$$\begin{aligned}\hat{c}_0 &= \sum_{m=0}^{\infty} \alpha^m c_m, \\ \hat{c}_1 &= (1 - \alpha^2) \sum_{m=1}^{\infty} m \alpha^{m-1} c_m, \\ \hat{c}_2 &= c_2 + \alpha(-c_1 + 3c_3) + \alpha^2(-4c_2 + 6c_4) + \cdots \alpha^3(c_1 - 9c_3 + 10c_5) + \cdots, \\ \hat{c}_3 &= c_3 + \alpha(-2c_2 + 4c_4) + \alpha^2(c_1 - 9c_3 + 10c_5) + \alpha^3(6c_2 - 24c_4 + 20c_6) + \cdots, \\ &\quad \cdots(3)\end{aligned}$$

$$\hat{c}_m = c_m + \begin{cases} (m+1)c_{m+1}\alpha, & m=0 \\ \{(m+1)c_{m+1} - (m-1)c_{m-1}\}\alpha, & m>0 \end{cases} \quad \cdots(4)$$

【0052】

式（4）は、式（2）で表される再帰式を、式（3）のように、伸縮パラメータ  $\alpha$  について展開し、 $\alpha$  の1次の項で近似し、上記文献4に記載されているように、HMMの最尤推定のためのQ関数に組み入れ、 $\alpha$  について、Q関数を最大にすることで導くことができる。

【0053】

導かれた関数が、式（5）である。

【 0 0 5 4 】

$$\alpha = \frac{\sum_{t=1}^T \gamma_t(j, k) \left[ \sum_{m=1}^M \frac{1}{\sigma_m^2} \Delta c_{mt} \bar{c}_{mt} \right]}{\sum_{t=1}^T \gamma_t(j, k) \left[ \sum_{m=1}^M \frac{1}{\sigma_m^2} \bar{c}_{mt}^2 \right]}$$

$$\Delta c_{mt} = c_{mt} - \mu_{jkm}, \quad \dots (5)$$

$$\bar{c}_{mt} = (m-1)c_{(m-1)t} - (m+1)c_{(m+1)t}$$

【 0 0 5 5 】

式 (5) において、 $c$  は、前記の入力パターンのうちのメルケプストラム部分を表し、 $\mu$  は、HMMの平均ベクトル表し、 $\sigma$  は、HMMの分散を表し、 $\gamma$  は、アライメント情報である時刻  $t$ 、状態  $j$ 、混合状態  $k$  における事後確率を表す。

【 0 0 5 6 】

事後確率は、フォワードバックワードアルゴリズムの場合、ある時刻状態における存在確率であり、ビタビアルゴリズムの場合、ある時刻状態に最適経路に存在している場合 1、存在していない場合 0 となる。

【 0 0 5 7 】

本実施例では、ワーピング関数として式 (1) をあげたが、本発明において、式 (1) に限定されるものでなく、任意の式が適用可能である。また、式 (5) を導くために、式 (2) の 1 次近似を用いたが、2 次以上の多次の近似を用いてもよい。

【 0 0 5 8 】

図 2 は、本発明の第 1 の実施例の処理を説明するための流れ図である。図 1 及び図 2 を参照して、本発明の第 1 の実施例の全体の動作について詳細に説明する。まず、音声信号と入力し (図 2 のステップ A 1 0 1)、分析部 1 で、入力パターンの計算を行う (ステップ A 1 0 2)。

【 0 0 5 9 】

伸縮推定部 3 が、分析部 1 から出力される入力パターンと、入力された HMM (A 1 0 5) とを用いて、伸縮パラメータを計算する (ステップ A 1 0 3)。

【 0 0 6 0 】

分析部 1 から出力された入力パターンを、変換部 2 が、式 (2)、式 (3)、

式(4)のいずれかの変換の関数を用い、変換入力パターンを求める(ステップA104)。ただし、 $\alpha$ の値は、初めての発声の場合0とし、2回目以降の発声の場合、伸縮推定部3が出力した値を用いる。

【0061】

次に本発明の第1の実施例の作用効果について説明する。本発明の第1の実施例では、分析部1の出力である入力パターンが、変換部2に入力されるように構成されているため、スペクトルの周波数の伸縮を、メルケプストラム領域で行うことができる。

【0062】

また、式(5)を用いた場合、従来の技術で説明したような、繰り返しの計算が無く、分析処理などが1回で済むため、伸縮パラメータを推定する際の演算量を少なくできる。

【0063】

次に、本発明の第2の実施例について説明する。図3は、本発明の第2の実施例の構成を示す図である。図3を参照すると、本発明の第2の実施例をなす音声認識装置は、分析部1と、変換部2と、伸縮推定部3と、認識部101と、標準パターン記憶部4と、を備えている。分析部1、変換部2、伸縮推定部3、標準パターン記憶部4は、前記第1の実施例で説明したものと同様とされる。

【0064】

すなわち、分析部1は、前記第1の実施例と同様、音声信号を分析し、入力パターンを計算し、出力する。

【0065】

変換部2は、前記第1の実施例と同様に、入力パターンの変換を行い、変換入力パターンとして出力する。

【0066】

標準パターン記憶部4には、前記第1の実施例と同様、音韻を表す要素として、入力パターンの平均ベクトルや分散等で構成されるHMMが記憶される。

【0067】

認識部(整合部)101は、HMMのうち、いずれが、変換部から出力された

変換入力パターンによく整合するかを調べることにより、認識を行い、認識結果を出力する。整合の方法として、上記文献 4 で示されるビタビアルゴリズムやフォワードバックワードアルゴリズムのような公知の方法が用いられる。

## 【 0 0 6 8 】

図 4 は、本発明の第 2 の実施例の処理手順を説明するための流れ図である。図 3 及び図 4 を参照して、本発明の第 2 の実施例の全体の動作について詳細に説明する。

## 【 0 0 6 9 】

分析部 1 が、入力された音声信号（図 4 のステップ B 1 0 1）を分析し、入力パターンの計算を行う（ステップ B 1 0 2）。分析部 1 から出力された入力パターンを、変換部 2 が、式（2）、式（3）、式（4）のいずれかの変換の関数を用い、変換入力パターンを求める（ステップ B 1 0 3）。ただし、 $\alpha$  の値は、初めての発声の場合 0 とし、2 回目以降の発声の場合、ワーピングパラメータの伸縮推定部 3 が出力した値を用いる。

## 【 0 0 7 0 】

つぎに、変換入力パターンを用い、認識部 1 0 1 が、認識処理を行う（ステップ B 1 0 4）。

## 【 0 0 7 1 】

このとき、認識部 1 0 1 に、標準パターン記憶部 4 から HMM が入力される（ステップ B 1 0 6）。認識処理後、認識結果を基に伸縮パラメータの推定部 0 0 3 が、伸縮パラメータを計算する（ステップ B 1 0 5）。

## 【 0 0 7 2 】

以後、ステップ B 1 0 5 で求めた伸縮パラメータを用い、ステップ B 1 0 1 の音声の入力の処理から繰り返す。

## 【 0 0 7 3 】

次に、本発明の第 2 の実施例の作用効果について説明する。本発明の第 2 の実施例では、前記第 1 の実施例のスペクトル変換部 1 0 0 と認識部 1 0 1 とを備えて構成されており、音声信号が入力されるたびに、伸縮パラメータの値が更新され、標準パターンとの周波数のずれを補正することができるため、認識性能を向

上している。

【 0 0 7 4 】

また本発明の第 2 の実施例では、伸縮パラメータの推定を HMM の最尤推定の Q 関数を最小にするような式 ( 5 ) を用いて行っているため、伸縮パラメータが連続的な値で推定され、予め用意された離散的な値を用いた場合に比べ、認識性能の向上を期待することができる。

【 0 0 7 5 】

次に、本発明の第 3 の実施例について説明する。図 5 は、本発明の第 3 の実施例の構成を示す図である。図 5 を参照すると、本発明の第 3 の実施例をなす標準パターン学習装置は、前記第 1 の実施例のスペクトル変換部 1 0 0 に加えて、学習音声記憶部 2 0 1 と、標準パターン推定部 2 0 2 と、尤度判定部 2 0 3 とを備えた構成されている。

【 0 0 7 6 】

学習用音声記憶部 2 0 1 は、HMM の学習に用いる音声信号を記憶する。

【 0 0 7 7 】

標準パターン推定部 2 0 2 は、スペクトル変換部 1 0 0 が出力する変換入力パターンと HMM を用い、HMM のパラメータを推定する。推定方法は、例えば、文献 4 等に記載されている、最尤推定を用いる。

【 0 0 7 8 】

尤度判定部 2 0 3 は、スペクトル変換部 1 0 0 が出力する変換入力パターンと HMM を用い、全ての学習用の音声信号に対し、距離を求める。距離を求める方法として、標準パターンが HMM なら、上記文献 5 に記載されているビタビアルゴリズムやフォワードバックワードアルゴリズムのような方法が用いられる。

【 0 0 7 9 】

本発明の第 3 の実施例では、HMM の学習を例として説明しているが、音声認識関係のパラメータの学習であれば、どのようなものでもよい。

【 0 0 8 0 】

図 6 は、本発明の第 3 の実施例の処理動作を説明するための流れ図である。図 5 及び図 6 を参照して、本発明の第 3 の実施例の全体の動作について詳細に説明



する。まず、学習用の音声信号を、スペクトル変換部 1 0 0 のスペクトル分析部 1 に入力する（図 6 の C 1 0 2）。

【 0 0 8 1 】

分析部 1 は、学習用の信号を分析し入力パターンを出力する（ステップ C 1 0 2）。伸縮推定部 0 0 3 が、伸縮パラメータを推定する（ステップ C 1 0 3）。

【 0 0 8 2 】

変換部 2 で、入力パターンの変換を行い、変換入力パターンを出力する（ステップ C 1 0 4）。

【 0 0 8 3 】

標準パターン推定部 2 0 2 において、変換入力パターンと HMM を用い、HMM の推定を行う（ステップ C 1 0 5）。

【 0 0 8 4 】

尤度判定部 2 0 3 が、全ての音声信号に対する尤度を求め、尤度の変化量とその閾値を比較し、尤度の変化量が閾値以下の場合、標準パターン推定部 2 0 2 で推定された HMM で標準パターン記憶部 4 を更新して、学習を終了する。

【 0 0 8 5 】

一方、尤度の変化量が、閾値より大きい場合、標準パターン推定部 2 0 2 で推定された HMM によって標準パターン記憶部 4 を更新し、学習音声データの入力処理（C 1 0 1）から、再び一連の処理を行う。

【 0 0 8 6 】

本発明の第 3 の実施例の作用効果について説明する。本発明の第 3 の実施例によれば、話者ごとの周波数の伸縮の影響を、ワーピング関数で修正した標準パターンを学習する際、伸縮パラメータの推定が学習処理の中でできるため、従来の技術よりも演算量を少なくできる。また、伸縮パラメータの推定に用いる式（5）は、HMM の最尤推定を用いて導出された式であるので、他の HMM のパラメータの推定と同様に、学習の過程に組み入れることが容易である。

【 0 0 8 7 】

次に本発明の第 4 の実施例について説明する。図 7 は、本発明の第 4 の実施例の構成を示す図である。図 7 を参照すると、本発明の第 4 の実施例は、前記第 1

の実施例の構成に加えて逆変換部 5 を備えている。逆変換部 5 は、変換部 2 から出力された伸縮後の前記入力パターンの時系列を逆変換して時間領域 (time domain) の信号波形を出力することで、声質の変換を行う。

## 【0088】

次に本発明の第 5 の実施例について説明する。図 8 は、本発明の第 5 の実施例の構成を示す図である。本発明の第 5 の実施例は、前記した第 1 乃至第 4 の実施例の装置を、コンピュータで実行されるプログラム制御によって実現するものである。図 8 を参照すると、図 1 に示した、分析部 1、変換部 2、伸縮推定部 3 の処理を、コンピュータ 10 上でプログラムを実行することで実現する場合、プログラムを記録した CD-ROM、DVD、FD、磁気テープ等の記録媒体 14 から記録媒体アクセス装置 13 を介してコンピュータ 10 の主記憶 12 にロードし、CPU 11 が該プログラムを実行する。すなわち、記録媒体 14 には、入力音声信号をケプストラムを含む入力パターンに変換する分析処理と、前記入力パターンと、標準パターンを記憶する標準パターン記憶部に記憶されている前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定処理と、前記伸縮パラメータを用いて前記入力パターンを変換する変換処理と、の各処理を前記コンピュータで実行させるためのプログラムが記録される。

## 【0089】

また、前記出力された伸縮後の前記入力パターンと前記標準パターンとの距離を算出し最も距離の小さい前記標準パターンを認識結果として出力する整合処理をコンピュータで実行させるためのプログラムを記録してもよい。

## 【0090】

さらに、記録媒体 14 には、学習用音声を記憶しておく学習音声記憶部に記憶されている学習用音声を入力しケプストラムを含む入力パターンに変換する分析処理と、前記入力パターンと、標準パターンを記憶する標準パターン記憶部に記憶されている前記標準パターンとを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定処理と、前記伸縮パラメータを用いて前記入力パターンを変換する変換処理と、学習用音声に対し、変換処理が出力した伸縮後の入力パターンと、標準パターンを利用して、前記標準パターンを更新する標準パターン推定処理と

、伸縮後の入力パターンと標準パターンを利用して、距離を算出し、距離の変化量を監視する尤度判定処理と、の各処理を前記コンピュータで実行させるためのプログラムを記録してもよい。このように、前記第 2 乃至第 4 の実施例についても、同様にして、プログラム制御で実現可能である。なお、プログラムは、不図示のサーバ装置からネットワーク等伝送媒体を介してダウンロードするようにしてもよい。すなわち、上記プログラムを担持する媒体であれば、記録媒体以外にも、通信媒体等任意の媒体を含む。

【0091】

【発明の効果】

以上説明したように、本発明によれば下記記載の効果を奏する。

【0092】

本発明の第 1 の効果は、音声信号のスペクトルの周波数の伸縮において、認識性能に最適なパラメータの計算に要する演算量を縮減する、ということである。

【0093】

その理由は、本発明においては、周波数軸に対し 1 次の全域通過型等のフィルタ処理における変換が、ケプストラム領域で伸縮パラメータべき級数の形で解け、その級数を 1 次関数と近似したとき、最尤推定のための Q 関数を最小にする伸縮パラメータの関数を容易な関数で記述することができ、該関数を用いて計算を行う構成としたためである。

【0094】

本発明の第 2 の効果は、HMM の学習時に、他のパラメータと同時に伸縮パラメータを推定できる、ということである。

【0095】

その理由は、本発明においては、前記第 1 の効果の理由で説明したように、伸縮パラメータを計算するための関数が、音声認識における最尤推定のための Q 関数から導出されるためである。

【図面の簡単な説明】

【図 1】

本発明の第 1 の実施例の構成を示す図である。

【図 2】

本発明の第 1 の実施例の処理動作を説明するための流れ図である。

【図 3】

本発明の第 2 の実施例の構成を示す図である。

【図 4】

本発明の第 2 の実施例の処理動作を説明するための流れ図である。

【図 5】

本発明の第 3 の実施例の構成を示す図である。

【図 6】

本発明の第 3 の実施例の処理動作を説明するための流れ図である。

【図 7】

本発明の第 4 の実施例の構成を示す図である。

【図 8】

本発明の第 5 の実施例の構成を示す図である。

【図 9】

従来のスペクトル変換部の構成を示す図である。

【図 1 0】

従来のスペクトル変換部の動作を示す流れ図である。

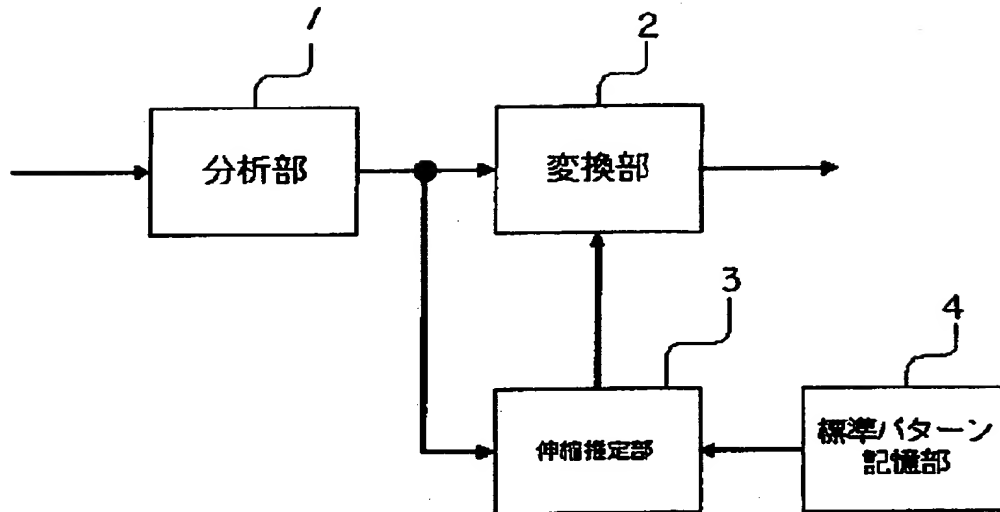
【符号の説明】

- 1 分析部
- 2 変換部
- 3 伸縮推定部
- 4 標準パターン記憶部
- 1 0 0 スペクトル変換部
- 1 0 1 認識部
- 2 0 1 学習音声記憶部
- 2 0 2 標準パターン推定部
- 2 0 3 尤度判定部
- 3 0 1 F F T 部

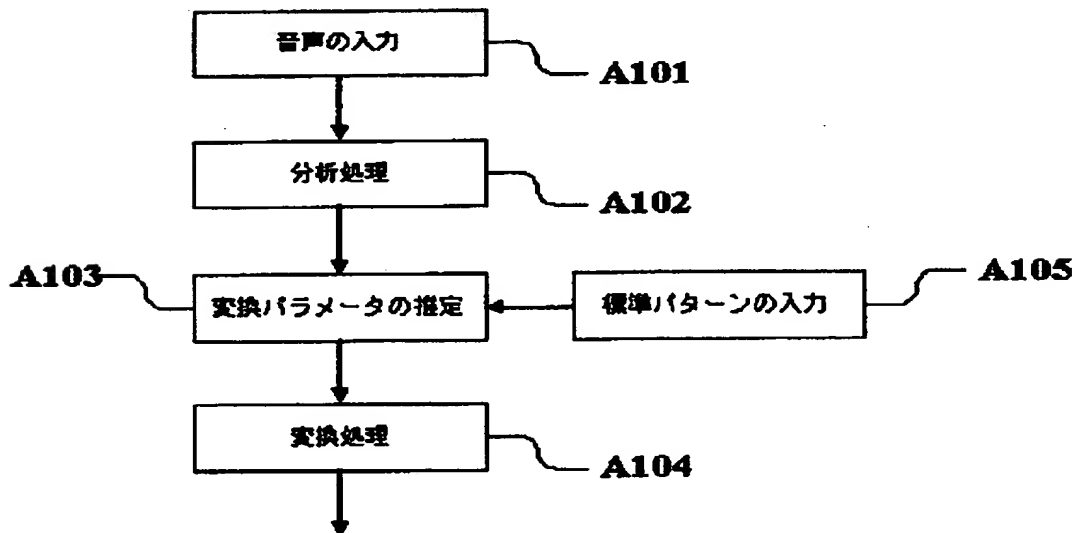
- 3 0 2 伸縮パラメータ記憶部
- 3 0 3 周波数変換部
- 3 0 4 特徴パラメータ計算部
- 3 0 5 標準パターン
- 3 0 6 整合部
- 3 0 7 伸縮パラメータ選択部

【書類名】 図面

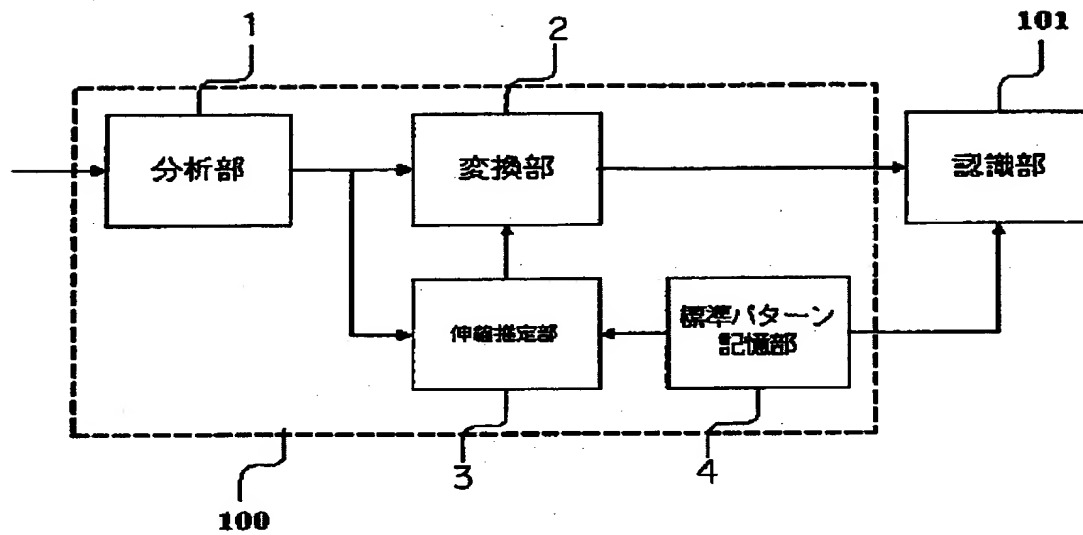
【図 1】



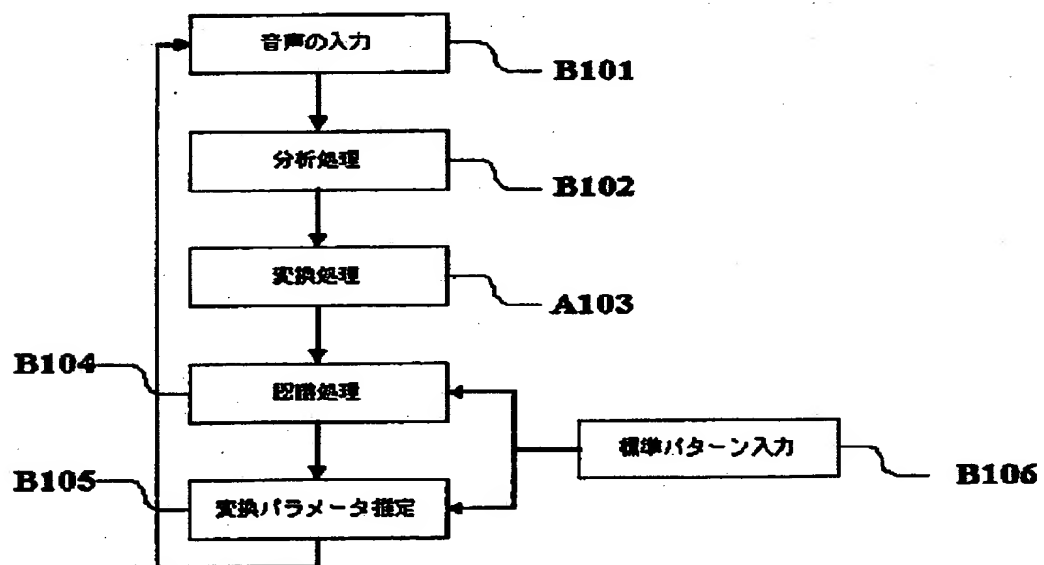
【図 2】



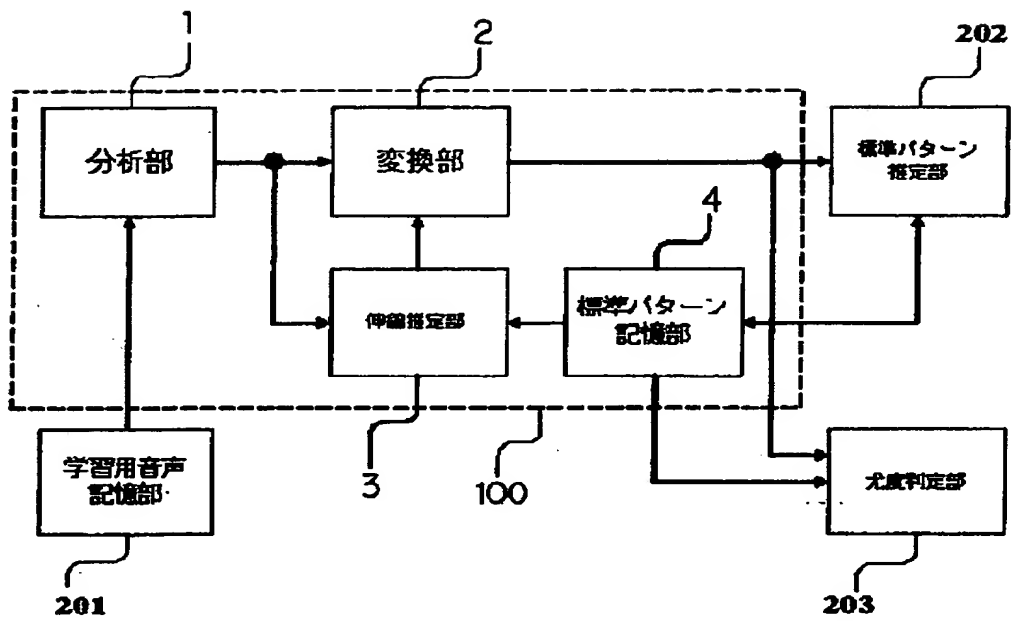
【図 3】



【図 4】

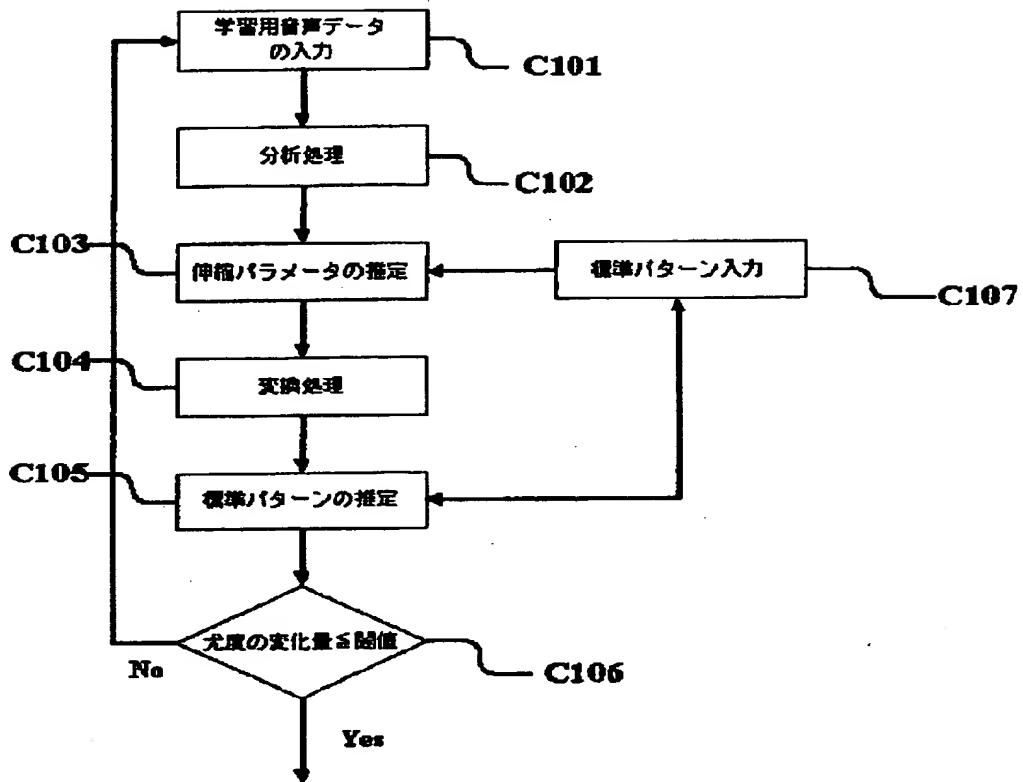


【図 5】

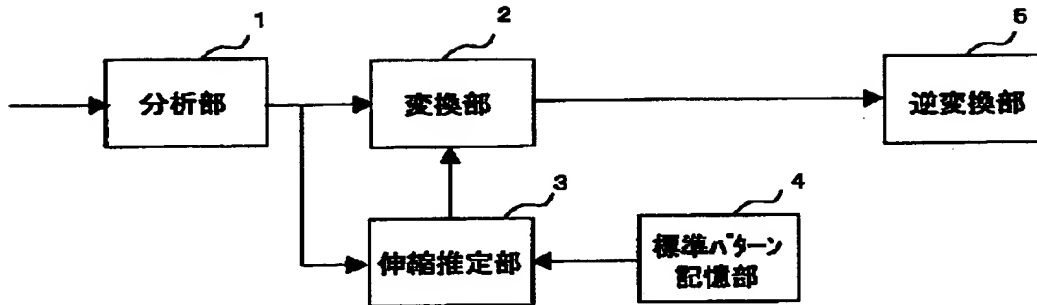




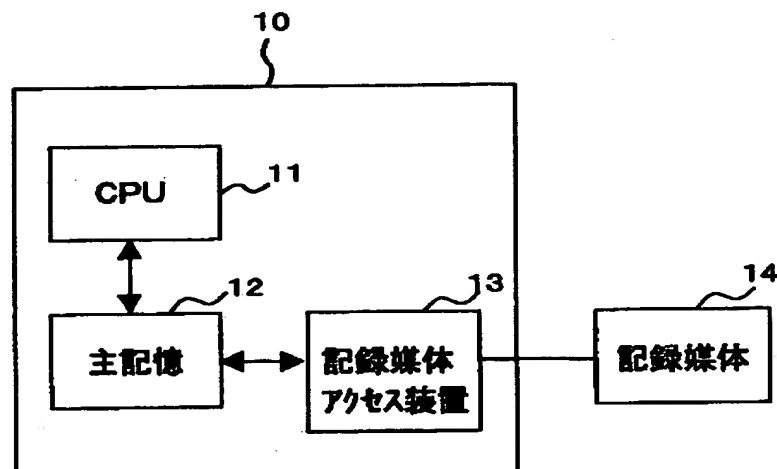
【図 6】



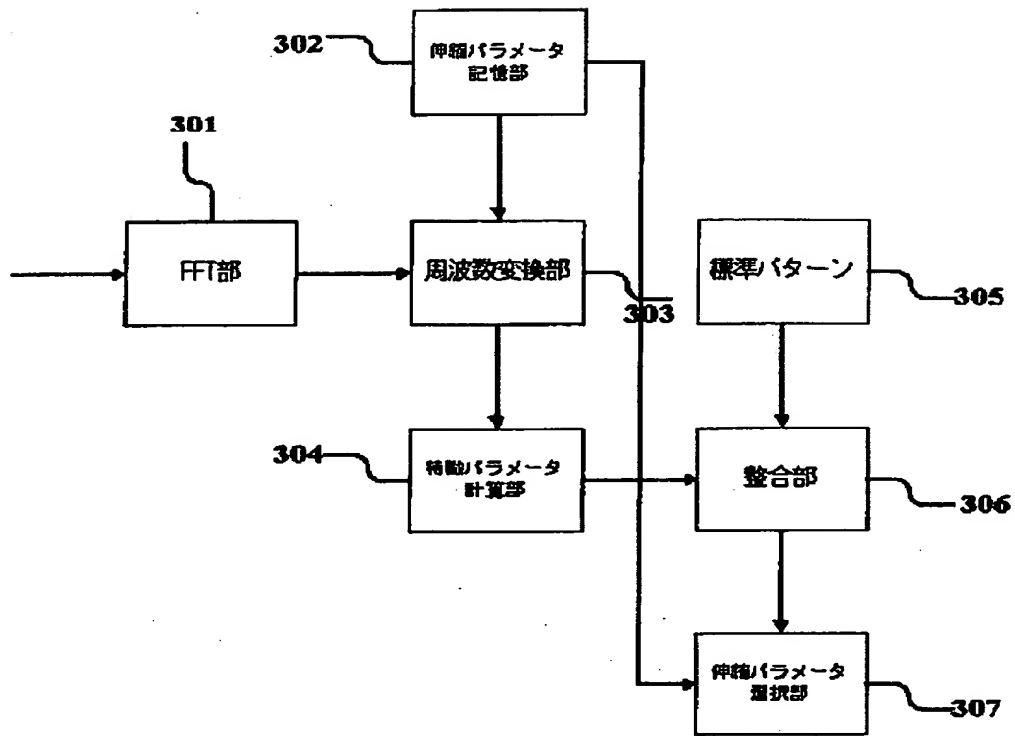
【図 7】



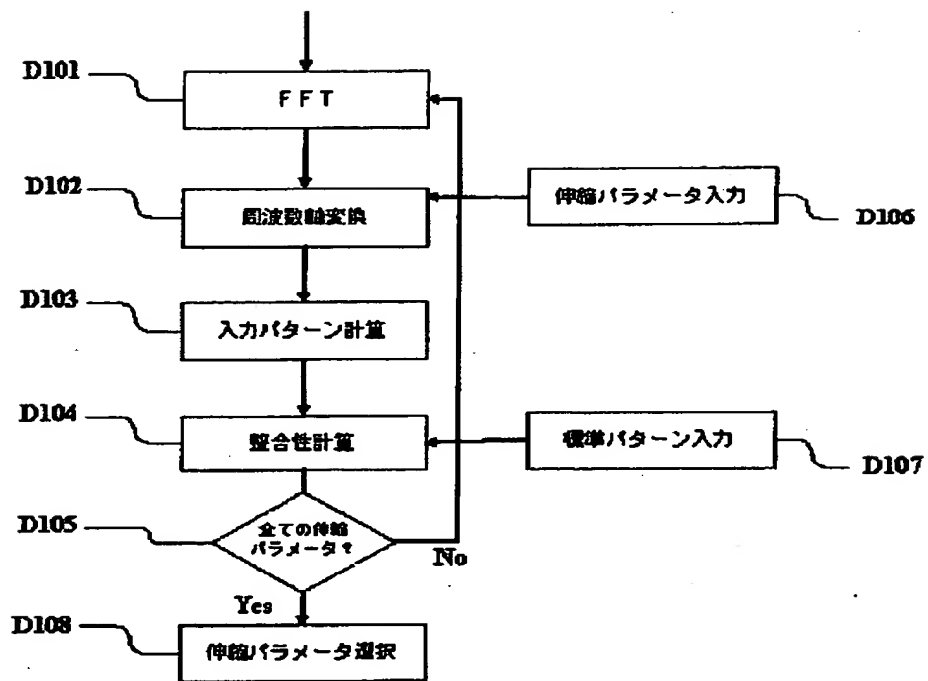
【図 8】



【図 9】



【図 1 0】



【書類名】 要約書

【要約】

【課題】

少ない演算量で話者毎に最適な伸縮パラメータの値を計算可能とし、性能を向上する装置及び方法の提供。

【解決手段】

入力音声信号をケプストラムを含む入力パターンに変換する分析部 1 と、標準パターンを記憶する標準パターン記憶部 3 と、入力パターンと標準パターンを用いて周波数軸方向の伸縮パラメータを出力する伸縮推定部 4 と、伸縮パラメータを用いて入力パターンを変換する変換部 2 と、変換部 2 によって変換された前記入力パターンと標準パターンとの距離を計算し最も距離の小さい標準パターンを認識結果として出力する認識部 1 0 1 を備え、伸縮推定部 4 は入力パターンに含まれるケプストラムを用いて伸縮パラメータを推定し、伸縮パラメータを決定する際に前もって様々な値を持つ必要が無く、様々な値について距離計算を行うことも不要とされる。

【選択図】

図 3

出 願 人 履 歴 情 報

識別番号 [0 0 0 0 0 4 2 3 7]

1. 変更年月日	1 9 9 0 年 8 月 2 9 日
[変更理由]	新規登録
住 所	東京都港区芝五丁目 7 番 1 号
氏 名	日本電気株式会社